

1. 49999 New York taxi trips



To drive a yellow New York taxi, you have to hold a "medallion" from the city's *Taxi and Limousine Commission*. Recently, one of those changed hands for over one million dollars, which shows how lucrative the job can be.

But this is the age of business intelligence and analytics! Even taxi drivers can stand to benefit from some careful investigation of the data, guiding them to maximize their profits. In this project, we will analyze a random sample of 49999 New York journeys made in 2013. We will also use regression trees and random forests to build a model that can predict the locations and times when the biggest fares can be earned.

Let's start by taking a look at the data!

```
In [2]: # Loading the tidyverse
library(tidyverse)

# Reading in the taxi data
taxi <- read_csv("datasets/taxi.csv")

head(taxi)

-- Attaching packages ----- tidyverse 1.
3.0 --
v ggplot2 3.3.2      v purrr   0.3.4
v tibble  3.0.3      v dplyr   1.0.2
v tidyr   1.1.2      v stringr 1.4.0
v readr   1.3.1      v forcats 0.5.0
-- Conflicts ----- tidyverse_conflict
s() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()
Parsed with column specification:
cols(
  medallion = col_character(),
  pickup_datetime = col_datetime(format = ""),
  pickup_longitude = col_double(),
  pickup_latitude = col_double(),
  trip_time_in_secs = col_double(),
  fare_amount = col_double(),
  tip_amount = col_double()
)

      medallion pickup_datetime pickup_longitude pickup_latitude tr
-----
4D24F4D8EF35878595044A52B098DFD2 2013-01-13
10:23:00 -73.94646 40.77273
A49C37EB966E7B05E69523D1CB7BE303 2013-01-13
04:52:00 -73.99827 40.74041
1E4B72A8E623888F53A9693C364AC05A 2013-01-13
10:47:00 -73.95346 40.77586
F7E4E9439C46B8AD5B16AB9F1B3279D7 2013-01-13
11:14:00 -73.98137 40.72473
A9DC75D59E0EA27E1ED328E8BE8CD828 2013-01-13
11:24:00 -73.96800 40.76000
19BF1BB516C4E992EA3FBAEDA73D6262 2013-01-13
10:51:00 -73.98502 40.76341
```

2. Cleaning the taxi data

As you can see above, the taxi dataset contains the times and price of a large number of taxi trips. Importantly we also get to know the location, the longitude and latitude, where the trip was started.

Cleaning data is a large part of any data scientist's daily work. It may not seem glamorous, but it makes the difference between a successful model and a failure. The taxi dataset needs a bit of polishing before we're ready to use it.

```
In [4]: # Renaming the location variables,
# dropping any journeys with zero fares and zero tips,
# and creating the total variable as the log sum of fare and tip
taxi <- taxi%>%
  rename(lat = pickup_latitude, long = pickup_longitude)%>%
  filter(fare_amount > 0 | tip_amount > 0)%>%
  mutate(total = log(fare_amount + tip_amount))

head(taxi)
```

	medallion	pickup_datetime	long	lat	trip_time_in_sec
	4D24F4D8EF35878595044A52B098DFD2	2013-01-13 10:23:00	-73.94646	40.77273	60
	A49C37EB966E7B05E69523D1CB7BE303	2013-01-13 04:52:00	-73.99827	40.74041	84
	1E4B72A8E623888F53A9693C364AC05A	2013-01-13 10:47:00	-73.95346	40.77586	6
	F7E4E9439C46B8AD5B16AB9F1B3279D7	2013-01-13 11:14:00	-73.98137	40.72473	72
	A9DC75D59E0EA27E1ED328E8BE8CD828	2013-01-13 11:24:00	-73.96800	40.76000	24
	19BF1BB516C4E992EA3FBAEDA73D6262	2013-01-13 10:51:00	-73.98502	40.76341	54

3. Zooming in on Manhattan

While the dataset contains taxi trips from all over New York City, the bulk of the trips are to and from Manhattan, so let's focus only on trips initiated there.

```
In [6]: # Reducing the data to taxi trips starting in Manhattan
# Manhattan is bounded by the rectangle with
# Latitude from 40.70 to 40.83 and
# Longitude from -74.025 to -73.93

taxi <- taxi %>%
  filter(lat >= 40.70, lat <= 40.83, long >= -74.025, long <= -73.93)
head(taxi)
```

	medallion	pickup_datetime	long	lat	trip_time_in_sec
	4D24F4D8EF35878595044A52B098DFD2	2013-01-13 10:23:00	-73.94646	40.77273	60
	A49C37EB966E7B05E69523D1CB7BE303	2013-01-13 04:52:00	-73.99827	40.74041	84
	1E4B72A8E623888F53A9693C364AC05A	2013-01-13 10:47:00	-73.95346	40.77586	6
	F7E4E9439C46B8AD5B16AB9F1B3279D7	2013-01-13 11:14:00	-73.98137	40.72473	72
	A9DC75D59E0EA27E1ED328E8BE8CD828	2013-01-13 11:24:00	-73.96800	40.76000	24
	19BF1BB516C4E992EA3FBAEDA73D6262	2013-01-13 10:51:00	-73.98502	40.76341	54

4. Where does the journey begin?

It's time to draw a map! We're going to use the excellent ggmap package together with ggplot2 to visualize where in Manhattan people tend to start their taxi journeys.

```
In [8]: # Loading in ggmap and viridis for nice colors
library(ggmap)
library(viridis)

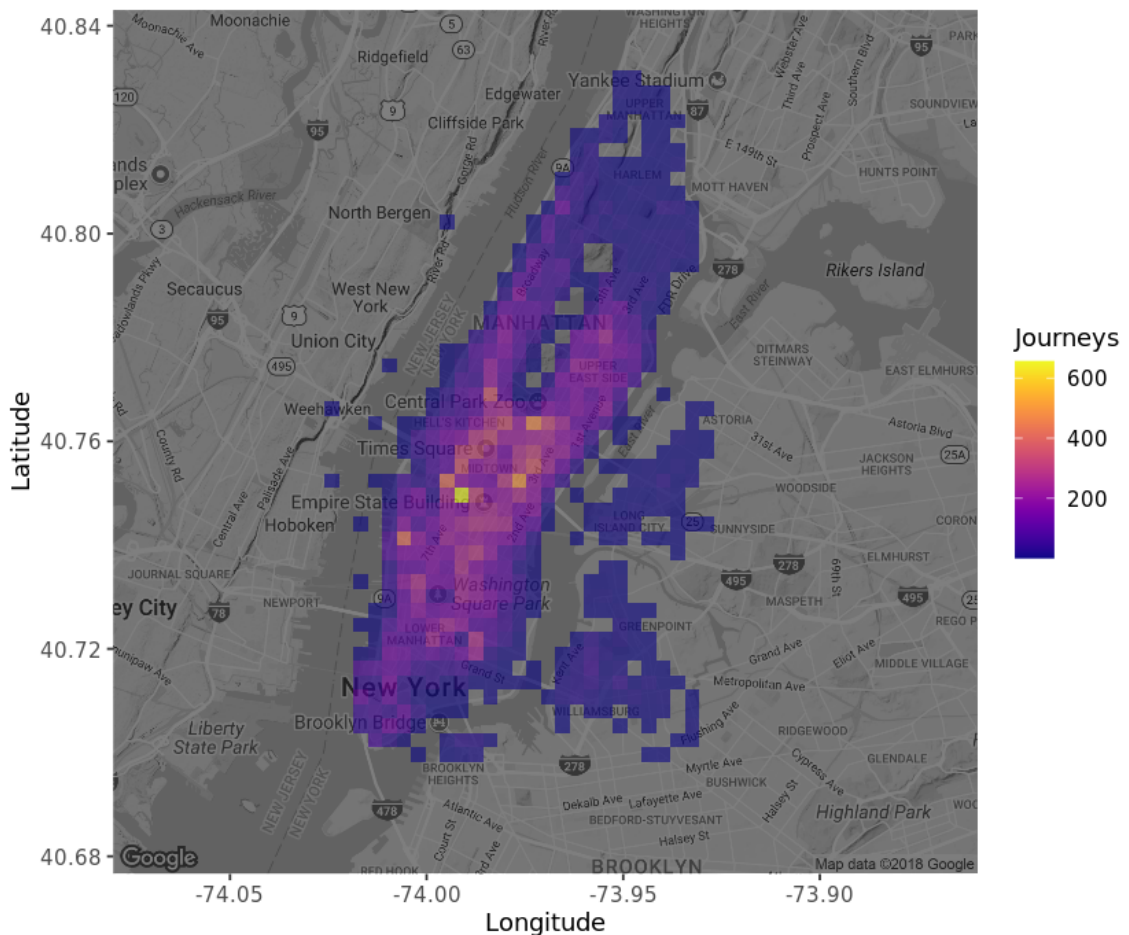
# Retrieving a stored map object which originally was created by
# manhattan <- get_map("manhattan", zoom = 12, color = "bw")
manhattan <- readRDS("datasets/manhattan.rds")

# Drawing a density map with the number of journey start locations
ggmap(manhattan, darken = 0.5) +
  scale_fill_viridis(option = 'plasma') +
  geom_bin2d(data = taxi, aes(long, lat), bins = 60, alpha = 0.6) +
  labs(x = 'Longitude', y = 'Latitude', fill = 'Journeys')
```

Google's Terms of Service: <https://cloud.google.com/maps-platform/terms/>.
(<https://cloud.google.com/maps-platform/terms/>.)

Please cite ggmap if you use it! See citation("ggmap") for details.

Loading required package: viridisLite



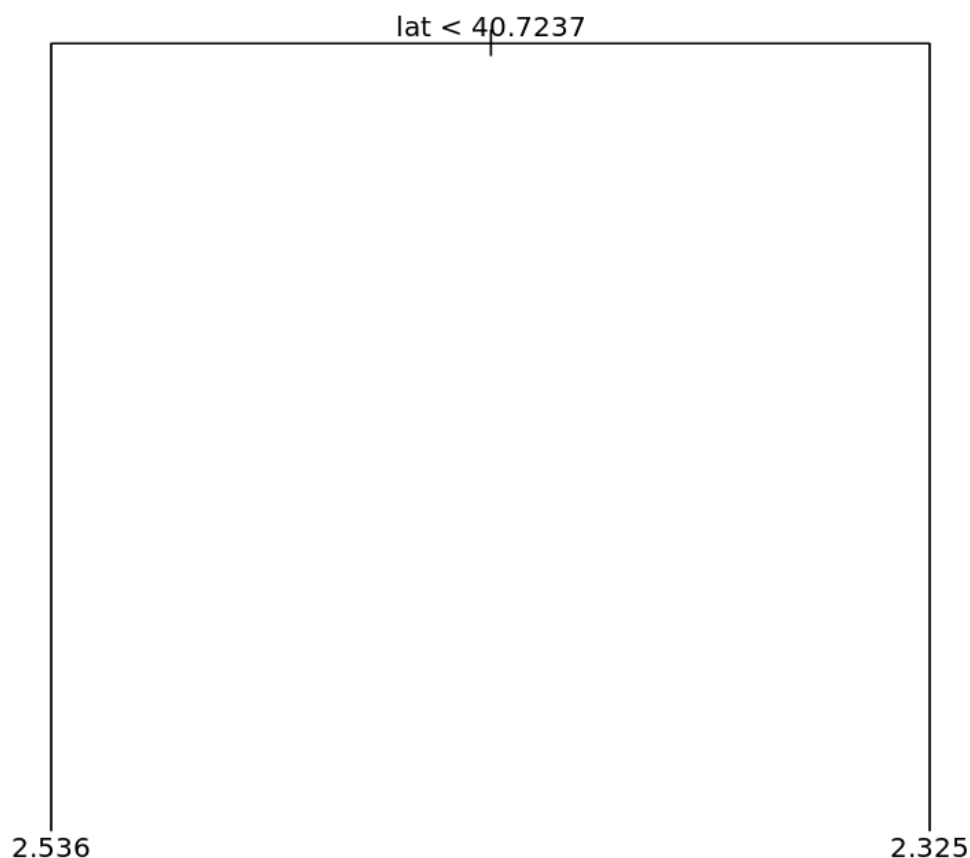
5. Predicting taxi fares using a tree

The map from the previous task showed that the journeys are highly concentrated in the business and tourist areas. We also see that some taxi trips originating in Brooklyn slipped through, but that's fine.

We're now going to use a regression tree to predict the total fare with lat and long being the predictors. The tree algorithm will try to find cutpoints in those predictors that results in the decision tree with the best predictive capability.

```
In [10]: # Loading in the tree package
library(tree)
# Fitting a tree to lat and long
fitted_tree <- tree(total ~ lat + long, data = taxi)

# Draw a diagram of the tree structure
plot(fitted_tree)
text(fitted_tree)
```



6. It's time. More predictors.

The tree above looks a bit frugal, it only includes one split: It predicts that trips where `lat < 40.7237` are more expensive, which makes sense as it is downtown Manhattan. But that's it. It didn't even include `long` as tree deemed that it didn't improve the predictions. Taxi drivers will need more information than this and any driver paying for your data-driven insights would be disappointed with that. As we know from Robert de Niro, it's best not to upset New York taxi drivers.

Let's start by adding some more predictors related to the *time* the taxi trip was made.

```
In [12]: # Loading in the lubridate package
library(lubridate)

# Generate the three new time variables
taxi <- taxi %>%
  mutate(hour = hour(pickup_datetime),
         wday = wday(pickup_datetime, label = TRUE),
         month = month(pickup_datetime, label = TRUE))
head(taxi)
```

Attaching package: 'lubridate'

The following object is masked from 'package:base':

date

	medallion	pickup_datetime	long	lat	trip_time_in_sec
	4D24F4D8EF35878595044A52B098DFD2	2013-01-13 10:23:00	-73.94646	40.77273	60
	A49C37EB966E7B05E69523D1CB7BE303	2013-01-13 04:52:00	-73.99827	40.74041	84
	1E4B72A8E623888F53A9693C364AC05A	2013-01-13 10:47:00	-73.95346	40.77586	6
	F7E4E9439C46B8AD5B16AB9F1B3279D7	2013-01-13 11:14:00	-73.98137	40.72473	72
	A9DC75D59E0EA27E1ED328E8BE8CD828	2013-01-13 11:24:00	-73.96800	40.76000	24
	19BF1BB516C4E992EA3FBAEDA73D6262	2013-01-13 10:51:00	-73.98502	40.76341	54

7. One more tree!

Let's try fitting a new regression tree where we include the new time variables.

```
In [14]: # Fitting a tree to lat and long

fitted_tree <- tree(total ~ lat + long + hour + wday + month, data = taxi)

# draw a diagram of the tree structure
plot(fitted_tree)
text(fitted_tree)

# Summarizing the performance of the tree
summary(fitted_tree)
```

Regression tree:

```
tree(formula = total ~ lat + long + hour + wday + month, data = taxi)
```

Variables actually used in tree construction:

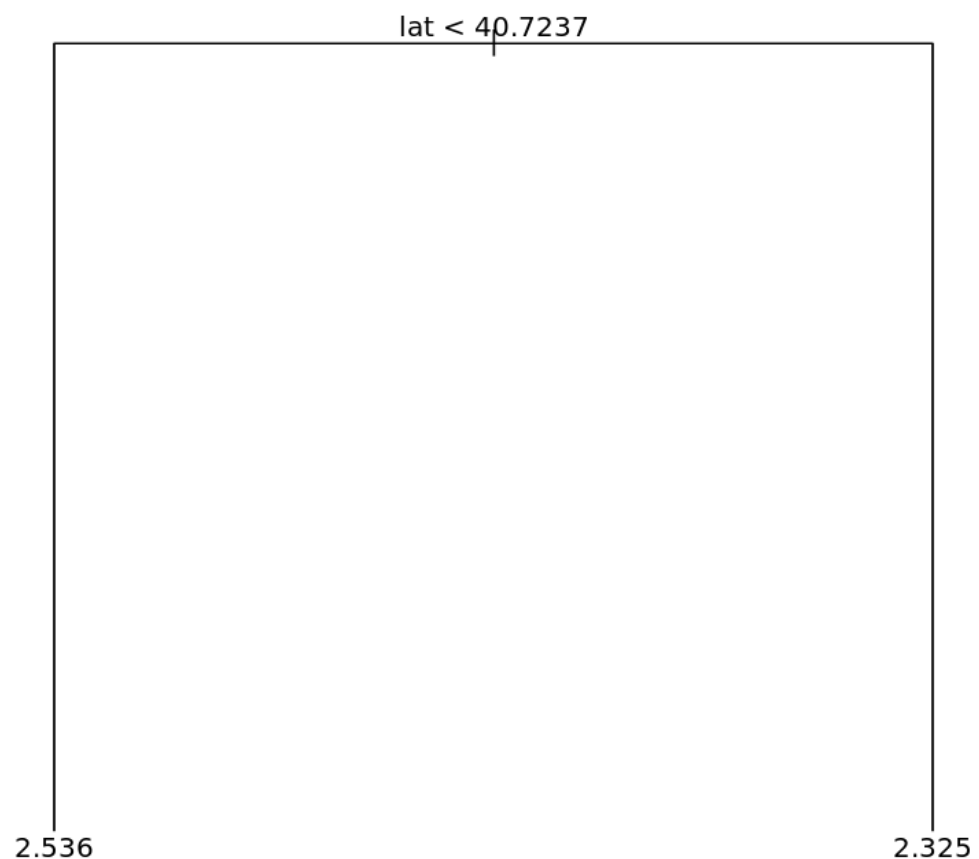
```
[1] "lat"
```

Number of terminal nodes: 2

Residual mean deviance: 0.3041 = 13910 / 45760

Distribution of residuals:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-1.61900	-0.37880	-0.04244	0.00000	0.32660	2.69900



8. One tree is not enough

The regression tree has not changed after including the three time variables. This is likely because latitude is still the most promising first variable to split the data on, and after that split, the other variables are not informative enough to be included. A random forest model, where many different trees are fitted to subsets of the data, may well include the other variables in some of the trees that make it up.

```
In [16]: # Loading in the randomForest package
library(randomForest)
# Fitting a random forest
fitted_forest <- randomForest(total ~ lat + long + hour + wday + month, data
                             ntree = 80,
                             sampsize = 10000)

# Printing the fitted_forest object
fitted_forest
```

randomForest 4.6-14

Type rfNews() to see new features/changes/bug fixes.

Attaching package: 'randomForest'

The following object is masked from 'package:dplyr':

combine

The following object is masked from 'package:ggplot2':

margin

Call:

```
randomForest(formula = total ~ lat + long + hour + wday + month,      dat
a = taxi, ntree = 80, sampsize = 10000)
```

Type of random forest: regression

Number of trees: 80

No. of variables tried at each split: 1

Mean of squared residuals: 0.3002836

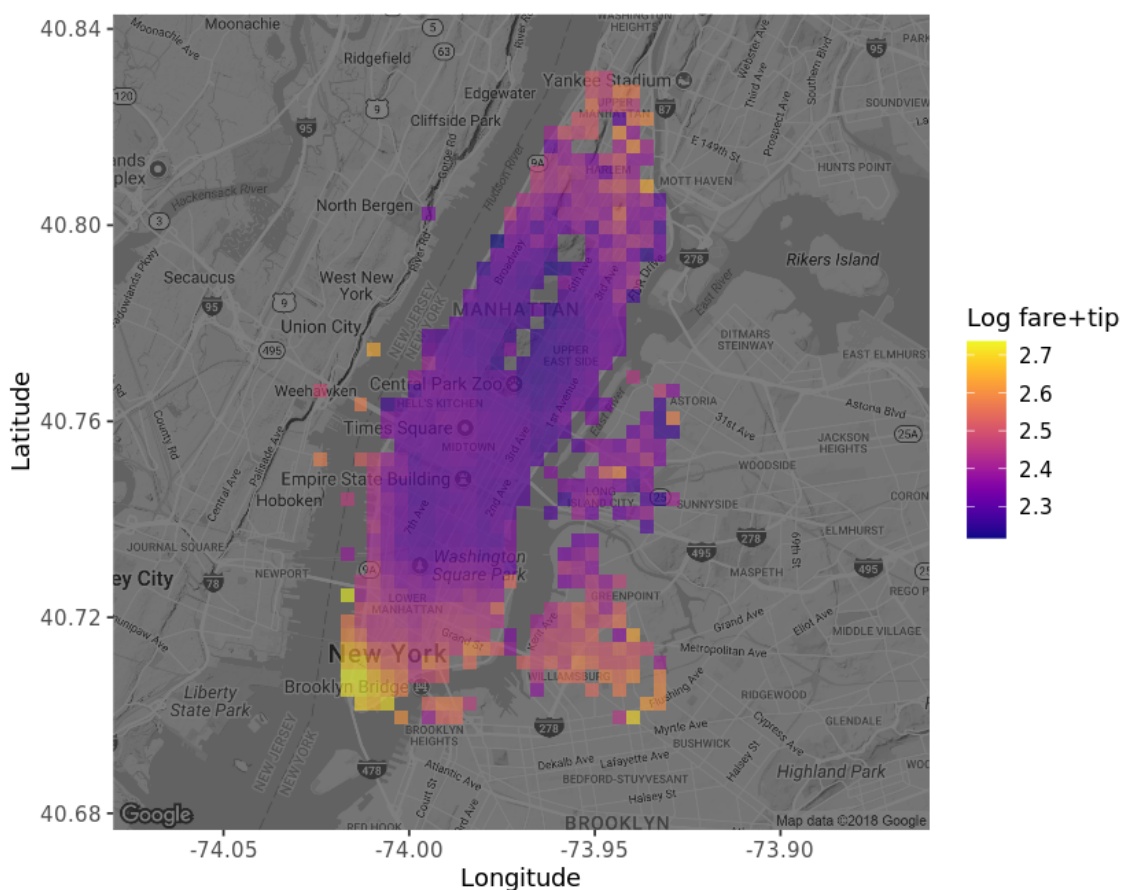
% Var explained: 2.61

9. Plotting the predicted fare

In the output of `fitted_forest` you should see the Mean of squared residuals, that is, the average of the squared errors the model makes. If you scroll up and check the summary of `fitted_tree` you'll find Residual mean deviance which is the same number. If you compare these numbers, you'll see that `fitted_forest` has a slightly lower error. Neither predictive model is *that* good, in statistical terms, they explain only about 3% of the variance.

Now, let's take a look at the predictions of `fitted_forest` projected back onto Manhattan.


```
In [18]: # Extracting the prediction from fitted_forest
taxi$pred_total <- fitted_forest$predicted
# Plotting the predicted mean trip prices from according to the random forest
ggmap(manhattan, darken = 0.5) +
  scale_fill_viridis(option = 'plasma') +
  stat_summary_2d(data = taxi, aes(x=long, y=lat, z=pred_total),
    fun = mean, bins = 60, alpha = 0.6) +
  labs(x = 'Longitude', y = 'Latitude', fill = 'Log fare+tip')
```



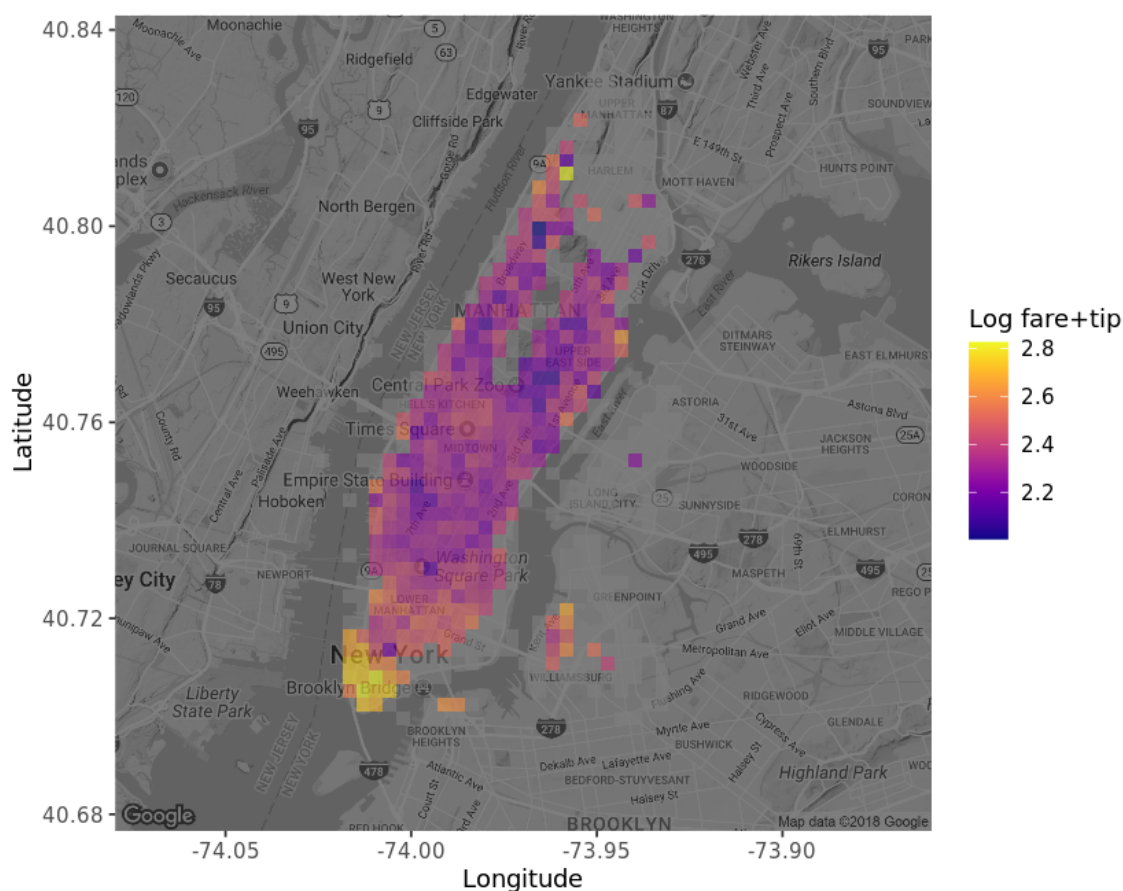
10. Plotting the actual fare

Looking at the map with the predicted fares we see that fares in downtown Manhattan are predicted to be high, while midtown is lower. This map only shows the prediction as a function of lat and long, but we could also plot the predictions over time, or a combination of time and space, but we'll leave that for another time.

For now, let's compare the map with the predicted fares with a new map showing the mean fares according to the data.

```
In [20]: # Function that returns the mean *if* there are 15 or more datapoints
mean_if_enough_data <- function(x) {
  ifelse( length(x) >= 15, mean(x), NA)
}

# Plotting the mean trip prices from the data
ggmap(manhattan, darken = 0.5) +
  scale_fill_viridis(option = 'plasma') +
  stat_summary_2d(data = taxi, aes(x=long, y=lat, z=total),
    fun = mean_if_enough_data, bins = 60, alpha = 0.6) +
  labs(x = 'Longitude', y = 'Latitude', fill = 'Log fare+tip')
```



11. Where do people spend the most?

So it looks like the random forest model captured some of the patterns in our data. At this point in the analysis, there are many more things we could do that we haven't done. We could add more predictors if we have the data. We could try to fine-tune the parameters of `randomForest`. And we should definitely test the model on a hold-out test dataset. But for now, let's be happy with what we have achieved!

So, if you are a taxi driver in NYC, where in Manhattan would you expect people to spend the most on a taxi ride?

```
In [22]: # Where are people spending the most on their taxi trips?
         spends_most_on_trips <- "downtown" # "uptown" or "downtown"
```